

Обучение машинного перевода без параллельных текстов*

Строганов А. А. Бахтеев¹ О. Ю. Стрижов² В. В.

¹Московский физико-технический институт

²Вычислительный центр им. А. А. Дородницына ФИЦ ИУ РАН

Данная работа посвящена исследованию машинного перевода без использования параллельных текстов. Исследование сконцентрировано на использовании нейросети с несколькими моделями Seq2seq для перевода с одного языка на другой и обратно. Особенностью данных моделей является то, что исходный текст кодируется во внутреннее представление модели, а затем декодируется в текст на другом языке. Две модели Seq2seq имеют общее скрытое пространство. Примером, иллюстрирующим работоспособность данного алгоритма, будет использован эксперимент по взаимному переводу с двух похожих языков – русского и украинского.

Ключевые слова: *машинный перевод, нейросеть, Seq2seq.*

1 Введение

Благодаря недавним достижениям в области глубокого обучения и наличию крупномасштабных параллельных корпусов, машинный перевод достиг впечатляющей производительности на нескольких языковых парах. Тем не менее, эти модели работают очень хорошо, только если они снабжены огромным количеством параллельных данных в порядке миллионов параллельных предложений. К сожалению, параллельные корпуса стоят дорого, поскольку они требуют специализированного опыта и часто не существуют для языков с низким уровнем ресурсов.

Есть несколько подходов к построению оптимального метода обучения. Предлагается использовать рекуррентные нейронные сети с короткой и долгой памятью и нейронные сети, в которых реализовано внимание. В других методах используются нейронные сети, которые осуществляют перевод в два этапа. Такой метода называется Seq2Seq.

Данная работа посвящена последнему методу последовательного перевода. Предлагается с помощью первой рекуррентной нейронной сети, основанной на долгой памяти перевести входящую последовательность в вектор, а с помощью второй перевести этот вектор в выходную последовательность на нужном нам языке. Данный метод позволяет гораздо быстрее обучить нейронную сеть переводу с одного языка на другой, в связи с использованием ей предыдущего опыта и наличию у нее памяти и внимания. Проверка и анализ метода проводятся с помощью алгоритма BLEU (Bilingual evaluation understudy) для проверки качества текста, переведенного с одного языка на другой на паре языков русский-украинский.

Литература

Поступила в редакцию

*Работа выполнена при финансовой поддержке РФФИ, проект №00-00-00000. Научный руководитель: Стрижов В. В. Авторы: А.В. Грабовой, О.Ю. Бахтеев, В.В. Стрижов, Eric Gaussier, координатор Малиновский Г.С. Консультант: Бахтеев О. Ю.